

CLAIMS:

1 1. A method for programming forwarding tables for switches for multipathing in a
2 subnet of a switched fabric including at least a host system, a target system and switches each
3 having one or more ports interconnected via links, said method comprising:

4 determining all possible links between all ports on the subnet during topology discovery;
5 creating an all port connectivity table which records all port-to-port connectivity
information;

6 creating an all switch shortest paths table which records all the shortest paths between
every switch pair on the subnet based the port-to-port connectivity information; and

7 computing forwarding tables for respective switches on the subnet that allow usage of
multiple paths between switch pairs based on the port-to-port connectivity information and based
on the shortest paths between every switch pair.

1 2. The method as claimed in claim 1, further comprising:

2 downloading the forwarding tables to respective switches on the subnet that allow usage
3 of multiple paths between switch pairs; and

4 enabling respective switches on the subnet to route data packets from the host system to
5 the target system via mutiple paths through the switched fabric.

1 3. The method as claimed in claim 1, wherein each of said host system and said
2 target system includes a channel adapter (CA) installed supporting one or more ports with each
3 port having multiple local identifiers (LIDs) assigned thereto for multipathing.

1 4. The method as claimed in claim 3, wherein each port on the subnet supports a
2 unique 16-bit LID and a LID Mask Control (LMC) which specifies the number of low order bits
3 of the LID to mask when checking a received destination LID against the port's destination LID.

1 5. The method as claimed in claim 1, wherein said forwarding tables are computed
2 to ensure loop-less paths and allow ports to be addressed by multiple local identifiers (LIDs), and
3 wherein said all-port connectivity and all-switch shortest paths tables are constantly updated
4 reflecting any dynamic changes to the subnet topology.

1 6. The method as claimed in claim 1, wherein said forwarding tables are computed
2 based on the principle that only the shortest path between a given switch pair is guaranteed to
3 overlap with other shortest paths that either originate from or destined to some intermediate port
4 that exists on the shortest path between the original switch pair.

1 7. The method as claimed in claim 1, wherein a forwarding table for a switch is
2 computed by:

1 determining a destination switch to which a destination port is directly connected,
2 identifying all the links that exist between the destination switch and other switches in the
3 subnet;
4 sorting all the links by respective originating port number in an ascending order;
5 picking an appropriate link and identifying the switch to which the link is connected at
6 the other end;
7 determining the best route between the switch identified and the switch for which the
8 forwarding table is being constructed; and
9 inputting associated outport number at a designated location in the forwarding table.

1 8. The method as claimed in claim 1, wherein said shortest paths between every
2 switch pair are computed utilizing an All Pair Shortest Paths (APSP) algorithm, and each shortest
3 path from the source to the destination switch is represented by a <Port, Cost> duple in which
4 port is the port number of the source switch where the path originates and cost is the path cost
5 metric that is computed based on a hop count, a message transfer (MTU) size, a link speed, width
6 and other port and link characteristics.

1 9. The method as claimed in claim 4, wherein a forwarding table for a switch is
2 computed by:
3 determining a destination switch to which a destination port is directly connected,

1 identifying all the links that exist between the destination switch and other switches in the
2 subnet;
3 sorting all the links by respective originating port number in an ascending order;
4 picking an appropriate link and identifying the switch to which the link is connected at
5 the other end;
6 determining the best route between the switch identified and the switch for which the
7 forwarding table is being constructed; and
8 inputting associated outport number at a designated location in the forwarding table.

10. The method as claimed in claim 4, wherein each of said forwarding tables for
1 switches in the subnet is computed, when multiple LIDs are assigned to channel adapter (CA)
2 ports, by a multipath assignment algorithm configured to:

3 receive information during the topology discovery including the number of multiple paths
4 (m) for configuration, the switch LID for which the forwarding table is being built (LID_s), and
5 the LID in the forwarding table which is currently identifying the correct outport (LID_d) for
6 forwarding data packets;

7 determine the base LID ($LID_{d\text{base}}$) for the given LID_d and determine if the base LID
8 ($LID_{d\text{base}}$) for the given LID_d corresponds to a switch using the all port connectivity table;
9 if the base LID ($LID_{d\text{base}}$) for the given LID_d corresponds to a switch using the all port
10 connectivity table, check if the base LID ($LID_{d\text{base}}$) for the given LID_d corresponds to the switch

- 1 LID for which the forwarding table is being built (LID_s);
 - 2 if the base LID ($LID_{d_{base}}$) for the given LID_d corresponds to the switch LID for which the
 - 3 forwarding table is being built (LID_s), indicate that there is no route for the given LID_d ;
 - 4 if the base LID ($LID_{d_{base}}$) for the given LID_d happens to be any switch other than the
 - 5 switch LID for which the forwarding table is being built (LID_s), check if the base LID ($LID_{d_{base}}$)
 - 6 for the given LID_d corresponds to the given LID_d ;
 - 7 if the base LID ($LID_{d_{base}}$) for the given LID_d corresponds to the given LID_d , set an arbitrary LID_x as the base LID ($LID_{d_{base}}$) for the given LID_d , and proceed to get the port number of switch LID_s for the best route between switch LID_s and LID_x from the all switch shortest paths table;
 - 8 if the base LID ($LID_{d_{base}}$) for the given LID_d does not correspond to a switch using the all port connectivity table, identify the peer node and port to which the port with $LID_{d_{base}}$ is connected from the all port connectivity table and determine if the peer node is a switch;
 - 9 if the peer node does not correspond to a switch, indicate that there is no route for the
 - 10 given LID_d ;
 - 11 if the peer node corresponds to a switch, determine if the peer switch LID ($LID_{s_{dest}}$)
 - 12 corresponds to the switch LID for which the forwarding table is being built (LID_s);
 - 13 if the peer switch LID ($LID_{s_{dest}}$) corresponds to LID_s , determine that the switch LID_s and
 - 14 port $LID_{d_{base}}$ are directly connected, and get the appropriate port number of switch LID_s from the
 - 15 all port connectivity table;

1 if the peer switch LID (LID_{sdest}) does not correspond to LID_s , identify all the links that are
2 directly connected to other switches from the peer switch LID (LID_{sdest}) and determine the
3 number of (N) links that connect switch LID_{sdest} where N is greater than "0";
4 if the number of (N) links is not greater than "0", indicate that there is no route for the
5 given LID_d ;
6 if the number of (N) links is greater than "0", set the offset (n) of LID_d from LID_{dbase} , and
7 determine if the number of multipaths (m) is less than the offset (n) of LID_d from LID_{dbase} and if
8 the number of (N) links that connect switch LID_{sdest} to other switch ports is less than the offset (n)
9 of LID_d from LID_{dbase} ;
10 if either the multipaths (m) or the (N) links is less than the offset (n) of LID_d from
11 LID_{dbase} , indicate that there is no route for the given LID_d ;
12 if neither the multipaths (m) nor the (N) links is less than the offset (n) of LID_d from
13 LID_{dbase} , store the LIDs of switches to which N links are connected in a list O(i) and sort the list
14 by the port number of switch LID_{sdest} from where each of N links originate in an ascending order;
15 set the list O(i).LID to the LID of the switch that the link O(n) is connected at the other
16 end, and check if O(i).LID corresponds to the switch LID for which the forwarding table is being
17 built (LID_s);
18 if the list O(i).LID corresponds to the switch LID for which the forwarding table is being
19 built (LID_s), set an arbitrary LIDx as the LID of the switch to which the port LID_d is directly
20 connected (LID_{sdest}) and then get the port number of switch LID_s for the best route between

1 switch LID_s and LIDx from the all switch shortest paths table;

2 if the list O(i).LID does not correspond to the switch LID for which the forwarding table

3 is being built (LID_s), set an arbitrary LIDx as the O(n).LID, and then get the port number of

4 switch LID_s for the best route between switch LID_s and LIDx from the all switch shortest paths

5 table; and

6 determine if the LID_d is the last LID assigned to a port or switch in the subnet, and

7 terminate computation of the forwarding table when LID_d is the last LID assigned to a port or

8 switch in the subnet.

1 11. A data network, comprising:

2 a host system having at least one channel adapter (CA) installed therein supporting one or

3 more ports with each port having multiple local identifiers (LIDs) assigned thereto for

4 multipathing;

5 at least one target system having at least one channel adapter (CA) installed therein

6 supporting one or more ports with each port having multiple local identifiers (LIDs) assigned

7 thereto for multipathing;

8 a switched fabric comprising a plurality of different switches which interconnect said host

9 system via CA ports to said remote system via CA port along different physical links for data

10 communications; and

11 a fabric manager provided in said host system for making topology discovery, assigning

1 local identifiers (LIDs) to all ports that are connected in the switched fabric, and programming
2 forwarding tables for switches in the switched fabric, wherein said fabric manager programs
3 forwarding tables for switches for multipathing by:

4 determining all possible links between all ports that are connected in the switched
5 fabric during topology discovery;
6 creating an all port connectivity table which records all port-to-port connectivity
7 information;

8 creating an all switch shortest paths table which records all the shortest paths
9 between every switch pair on the switched fabric based the port-to-port connectivity
10 information; and

11 computing forwarding tables for respective switches on the switched fabric that
12 allow usage of multiple paths between switch pairs based on the port-to-port connectivity
13 information and based on the shortest paths between every switch pair.

12. The data network as claimed in claim 11, wherein said fabric manager is
13 configured to download the forwarding tables to respective switches on the switched fabric that
14 allow usage of multiple paths between every switch pair, and enable respective switches on the
15 switched fabric to route data packets from the host system to the target system via multiple paths
16 through the switched fabric.

1 13. The data network as claimed in claim 11, wherein each port on the subnet
2 supports a unique 16-bit LID and a LID Mask Control (LMC) which specifies the number of low
3 order bits of the LID to mask when checking a received destination LID against the port's
4 destination LID.

1 14. The data network as claimed in claim 11, wherein said forwarding tables are
2 computed to ensure loop-less paths and allow ports to be addressed by multiple local identifiers
3 (LIDs), and wherein said all-port connectivity and all-switch shortest paths tables are constantly
4 updated reflecting any dynamic changes to the subnet topology.

1 15. The data network as claimed in claim 11, wherein said forwarding tables are
2 computed based on the principle that only the shortest path between a given switch pair is
3 guaranteed to overlap with other shortest paths that either originate from or destined to some
4 intermediate port that exists on the shortest path between the original switch pair.

1 16. The data network as claimed in claim 11, wherein said fabric manager is
2 configured to compute a forwarding table for a switch by:
3 determining a destination switch to which a destination port is directly connected,
4 identifying all the links that exist between the destination switch and other switches in the
5 subnet;

- sorting all the links by respective originating port number in an ascending order;
- picking an appropriate link and identifying the switch to which the link is connected at one end;
- determining the best route between the switch identified and the switch for which the forwarding table is being constructed; and
- inputting associated outport number at a designated location in the forwarding table.

17. The data network as claimed in claim 11, wherein said shortest paths between every port pair are computed utilizing an All Pair Shortest Paths (APSP) algorithm, and each shortest path from the source to the destination switch is represented by a <Port, Cost> duple in which port is the port number of the source switch where the path originates and cost is the path cost metric that is computed based on a hop count, a message transfer (MTU) size, a link speed, width and other port and link characteristics.

18. The data network as claimed in claim 11, wherein said fabric manager is provided with an algorithm for computing forwarding tables for switches, when multiple LIDs are assigned to channel adapter (CA) ports, by:

receiving information during the topology discovery including the number of multiple paths (m) for configuration, the switch LID for which the forwarding table is being built (LID_s), and the LID in the forwarding table which is currently identifying the correct outport (LID_d) for

1 forwarding data packets;

2 determining the base LID ($LID_{d_{base}}$) for the given LID_d and if the base LID ($LID_{d_{base}}$) for
3 the given LID_d corresponds to a switch using the all port connectivity table;

4 if the base LID ($LID_{d_{base}}$) for the given LID_d corresponds to a switch using the all port
5 connectivity table, checking if the base LID ($LID_{d_{base}}$) for the given LID_d corresponds to the
6 switch LID for which the forwarding table is being built (LID_s);

7 if the base LID ($LID_{d_{base}}$) for the given LID_d corresponds to the switch LID for which the
8 forwarding table is being built (LID_s), indicating that there is no route for the given LID_d ;

9 if the base LID ($LID_{d_{base}}$) for the given LID_d happens to be any switch other than the
10 switch LID for which the forwarding table is being built (LID_s), checking if the base LID
11 ($LID_{d_{base}}$) for the given LID_d corresponds to the given LID_d ;

12 if the base LID ($LID_{d_{base}}$) for the given LID_d corresponds to the given LID_d , setting an
13 arbitrary LID_x as the base LID ($LID_{d_{base}}$) for the given LID_d , and getting the port number of
14 switch LID_s for the best route between switch LID_s and LID_x from the all switch shortest paths
15 table;

16 if the base LID ($LID_{d_{base}}$) for the given LID_d does not correspond to a switch using the all
17 port connectivity table, identifying the peer node and port to which the port with $LID_{d_{base}}$ is
18 connected from the all port connectivity table and determining if the peer node is a switch;

19 if the peer node does not correspond to a switch, indicating that there is no route for the
20 given LID_d ;

1 if the peer node corresponds to a switch, determining if the peer switch LID (LID_{sdest})
2 corresponds to the switch LID for which the forwarding table is being built (LID_s);
3 if the peer switch LID (LID_{sdest}) corresponds to LID_s , determining that the switch LID_s
4 and port LID_{base} are directly connected, and getting the appropriate port number of switch LID_s
5 from the all port connectivity table;
6 if the peer switch LID (LID_{sdest}) does not correspond to LID_s , identifying all the links that
7 are directly connected to other switches from the peer switch LID (LID_{sdest}) and determining the
8 number of (N) links that connect switch LID_{sdest} where N is greater than "0";
9 if the number of (N) links is not greater than "0", indicating that there is no route for the
10 given LID_d;
11 if the number of (N) links is greater than "0", setting the offset (n) of LID_d from LID_{base},
12 and determining if the number of multipaths (m) is less than the offset (n) of LID_d from LID_{base}
13 and if the number of (N) links that connect switch LID_{sdest} to other switch ports is less than the
14 offset (n) of LID_d from LID_{base};
15 if either the multipaths (m) or the (N) links is less than the offset (n) of LID_d from
16 LID_{base}, indicating that there is no route for the given LID_d;
17 if neither the multipaths (m) nor the (N) links is less than the offset (n) of LID_d from
18 LID_{base}, storing the LIDs of switches to which N links are connected in a list O(i) and sorting the
19 list by the port number of switch LID_{sdest} from where each of N links originate in an ascending
20 order;

1 setting the list $O(i).LID$ to the LID of the switch that the link $O(n)$ is connected at the
2 other end, and checking if $O(i).LID$ corresponds to the switch LID for which the forwarding table
3 is being built (LID_s);

4 if the list $O(i).LID$ corresponds to the switch LID for which the forwarding table is being
5 built (LID_s), setting an arbitrary LID_x as the LID of the switch to which the port LID_d is directly
6 connected (LID_{sdest}) and then getting the port number of switch LID_s for the best route between
7 switch LID_s and LID_x from the all switch shortest paths table;

8 if the list $O(i).LID$ does not correspond to the switch LID for which the forwarding table
is being built (LID_s), setting an arbitrary LID_x as the $O(n).LID$, and then getting the port number
of switch LID_s for the best route between switch LID_s and LID_x from the all switch shortest
paths table; and

9 determining if the LID_d is the last LID assigned to a port or switch in the subnet, and
terminating computation of the forwarding table when LID_d is the last LID assigned to a port or
switch in the subnet.

19. A computer readable medium comprising instructions that, when executed by a
computer system, cause the computer system to:

3 determine all possible links between all ports on a subnet including at least a host system,
4 a target system and switches each having one or more ports interconnected via links during
5 topology discovery;

1 create an all port connectivity table which records all port-to-port connectivity
2 information;
3 create an all switch shortest paths table which records all the shortest paths between every
4 port pair on the subnet based the port-to-port connectivity information; and
5 compute forwarding tables for respective switches on the subnet that allow usage of
6 multiple paths between port pairs based on the port-to-port connectivity information and based
7 on the shortest paths between every port pairs.

1 20. The computer readable medium as claimed in claim 19, further causing the
2 computer system to:
3 download the forwarding tables to respective switches on the subnet that allow usage of
4 multiple paths between port pairs; and
5 enable respective switches on the subnet to route data packets from the host system to the
6 target system via multiple paths through the switched fabric.

1 21. The computer readable medium as claimed in claim 19, wherein each port on the
2 subnet supports a unique 16-bit LID and a LID Mask Control (LMC) which specifies the number
3 of low order bits of the LID to mask when checking a received destination LID against the port's
4 destination LID.

1 22. The computer readable medium as claimed in claim 19, wherein said forwarding
2 tables are computed to ensure loop-less paths and allow ports to be addressed by multiple local
3 identifiers (LIDs), and wherein said all-port connectivity and all-switch shortest paths tables are
4 constantly updated reflecting any dynamic changes to the subnet topology.

1 23. The computer readable medium as claimed in claim 19, wherein said forwarding
2 tables are computed based on the principle that only the shortest path between a given port pair is
3 guaranteed to overlap with other shortest paths that either originate from or destined to some
4 intermediate port that exists on the shortest path between the original port pair.

5 24. The computer readable medium as claimed in claim 19, wherein a forwarding
6 table for a switch is computed by:
7 determining a destination switch to which a destination port is directly connected,
8 identifying all the links that exist between the destination switch and other switches in the
9 subnet;
10 sorting all the links by respective originating port number in an ascending order;
11 picking an appropriate link and identifying the switch to which the link is connected at
12 the other end;
13 determining the best route between the switch identified and the switch for which the
14 forwarding table is being constructed; and

1 inputting associated outport number at a designated location in the forwarding table.

1
2 25. The computer readable medium as claimed in claim 19, wherein said shortest
3 paths between every switch pair are computed utilizing an All Pair Shortest Paths (APSP)
4 algorithm, and each shortest path from the source to the destination switch is represented by a
5 <Port, Cost> duple in which port is the port number of the source switch where the path
6 originates and cost is the path cost metric that is computed based on a hop count, a message
transfer (MTU) size, a link speed, width and other port and link characteristics.